



White Paper
Server Platform Group
May 2008, Intel Corporation

Intel® QuickData Technology Software Guide for Linux*

Overview

Contents

1. Overview	2
2. Intel® QuickData Technology	2
2.1 Faster, Scalable, More Reliable I/O	2
3. Supported Platforms	2
4. Requirements to Enable Intel QuickData Technology	2
4.1 How to Enable Intel QuickData Technology?	3
5. Intel QuickData Technology Module Overview	3
6. Intel QuickData Technology Client Interfaces	3
6.1 Client Registration/Un-registration	3
6.2 Initiating DMA Channel memcpy Operations	4
6.3 Client Data Structures	6
7. Intel QuickData Technology sysfs Interface	7
8. Support	7
9. References	7

This document describes the Intel® QuickData Technology and contains instructions on how to configure the Intel QuickData Technology for Linux* operating system.

2. Intel® QuickData Technology

Intel QuickData Technology is a platform solution designed to maximize the throughput of server data traffic across a broader range of configurations and server environments to achieve faster, scalable, and more reliable I/O.

2.1 Faster, Scalable, More Reliable I/O

Intel QuickData Technology is a component of Intel® I/O Acceleration Technology (Intel® I/OAT). It provides fast, scalable, and reliable throughput by moving data more efficiently through the server. By enabling broader use of the acceleration engine, Intel is improving the cost/performance of Intel® Xeon® 5000 series processor-based server platforms.

Intel has extended the value of Intel QuickData Technology by enabling networking solutions providers to utilize Intel QuickData Technology, the data movement acceleration engine present in Intel Xeon processor-based server platforms.

3. Supported Platforms

Intel QuickData Technology is supported on the following chipsets:

- Intel® 5000 Series Chipset
- Intel® 5100 Chipset
- Intel® 7300 Chipset
- Intel® 5400 Chipset

4. Requirements to Enable Intel QuickData Technology

This section describes the basic platform requirement to enable Intel QuickData Technology for Linux operating system.

1. Platform BIOS shall support Intel QuickData Technology.
2. Linux kernel 2.6.18 or above is required.
3. Make sure that kernel is build with Intel I/OAT DMA support. Kernel can be configured to build the IOATDMA as a module or as in-built driver.
4. To set TCP stack as Intel I/OAT DMA client, TCP receive copy offload should be enabled in kernel configuration.

4.1 How to Enable Intel QuickData Technology?

The Intel QuickData Technology can be enabled by loading the Intel I/OAT DMA engine driver/modules (ioatdma.ko) at runtime.

To enable Intel QuickData Technology, load the ioatdma driver module by typing following at the prompt.

```
modprobe ioatdma
```

Removing the ioatdma module once it has been loaded is not recommended since TCP receives a hold reference to the ioatdma driver when offloading traffic. However, if the "Forced module unloading" option is enabled in the kernel, the module may be unloaded with:

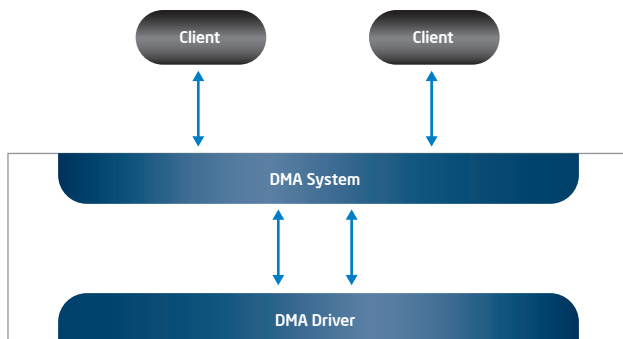
```
rmmod -f ioatdma
```

WARNING: The command above may hang the system.

5. Intel QuickData Technology Module Overview

This section highlights the Linux Intel QuickData Technology module (ioatdma.ko) implementation details.

Linux Intel QuickData Technology Module (IOATDMA)



1. DMA Driver

The DMA driver handles the hardware-specific details of performing DMA transactions. The DMA driver hides all the minute hardware details from upper layers, and handles the detection of Intel QuickData Technology, enumeration of resources, and handling of interrupts and power management events.

2. DMA System

The DMA subsystem abstracts the hardware-specific details from Intel QuickData Technology clients, and provides a registration interface. It manages all the underlying protocol to handle DMA channels.

3. Client

DMA clients register with the DMA subsystem before the DMA driver has been loaded.

6. Intel QuickData Technology Client Interfaces

This section defines the interfaces and data structures used by Intel QuickData Technology clients.

6.1 Client Registration/Un-registration

dma_async_client_register

This function registers Intel QuickData Technology client's function callback with the IOATDMA module.

void

```
dma_async_client_register(struct dma_client *client)
```

Parameter Description	
client	The pointer to a dma_client structure with event callback pointer and capability mask.

Return Values:

- None

dma_async_client_chan_request

This function is used by Intel QuickData Technology clients to request one or more DMA channels.

void

```
dma_async_client_chan_request(struct dma_client *client)
```

Parameter Description	
client	Intel QuickData Technology client identifier, passed to dma_async_client_register() (see above).

Return Values:

- None

dma_async_client_unregister

This function is used by Intel QuickData Technology client to unregister and to free all the resources in use.

void

dma_async_client_unregister(struct dma_client *client)

Parameter Description	
client	The pointer to a struct dma_client, passed to dma_async_client_register.

Return Values:

- None

6.2 Initiating DMA Channel memcpy Operations

Once the client has been given a DMA channel, it may use the following functions to perform accelerated memory copy operations.

dma_async_memcpy_buf_to_buf

This function initiates a buffer to buffer memory copy.

dma_cookie_t

dma_async_memcpy_buf_to_buf(struct dma_chan *chan, void *dest, void *src, size_t len)

Parameter Description	
chan	An allocated dma channel pointer.
dest	The destination address for the operation. This is a kernel virtual address. If the operation is to user memory, the caller must pin the memory, and pass a kernel pointer.
src	The source address for the operation. This is a kernel virtual address. If the operation is from user memory, the caller must pin the memory, and pass a kernel pointer.
len	The length of the memcpy.

Return Values:

- If negative, the return is a standard Linux kernel error code.
- If non-negative, the return is a "cookie". Since memory copies happen asynchronously, the caller can use this function to verify when a transaction, and all transactions that preceded it, is complete. See dma_async_memcpy_complete.

dma_async_memcpy_buf_to_pg

This function initiates a buffer to page memory copy.

dma_cookie_t

dma_async_memcpy_buf_to_pg(struct dma_chan *chan, struct page *page, unsigned int offset, void *kdata, size_t len)

Parameter Description	
chan	An allocated dma channel pointer.
page	The destination page for the operation. This is a kernel virtual address. If the operation is to user memory, the caller must pin the memory, and pass a kernel pointer.
offset	Offset into the page where to put the data.
kdata	The source address for the operation. This is a kernel virtual address.
len	The length of the memcpy.

Return Values:

- If negative, the return is a standard Linux kernel error code.
- If non-negative, the return is a "cookie". Since memory copies happen asynchronously, the caller can use this function to verify when a transaction, and all transactions that preceded it, is complete. See dma_async_memcpy_complete.

dma_async_memcpy_pg_to_pg

This function initiates a page to page memory copy.

dma_cookie_t

```
dma_async_memcpy_pg_to_pg(struct dma_chan *chan,
    struct page *dest_pg, unsigned int dest_off,
    struct page *src_pg, unsigned int src_off, size_t len)
```

Parameter Description	
chan	An allocated dma channel pointer.
dest_pg	The destination page for the operation, must be mappable according to DMA mapping API rules for streaming mappings.
dest_off	Offset into the page where to put the data.
src_pg	The source page for the operation, must be mappable according to DMA mapping API rules for streaming mappings.
src_off	Offset into the page where to find the data.
len	The length of the memcpy.

Return Values:

- If negative, the return is a standard Linux kernel error code.
- If non-negative, the return is a "cookie". Since memory copies happen asynchronously, the caller can use this function to verify when a transaction, and all transactions that preceded it, is complete. See dma_memcpy_complete.

dma_async_memcpy_issue_pending

This function flushes any pending copy requests to the hardware copy device. This allows drivers to push copies to hardware in batches, reducing MMIO writes where possible.

```
void dma_async_memcpy_issue_pending(struct dma_chan *chan)
```

Parameter Description	
chan	An allocated dma channel pointer.

dma_async_is_complete

This function tests a DMA request for completion.

enum dma_status

```
dma_async_is_complete(dma_cookie_t cookie,
    dma_cookie_t last_complete, dma_cookie_t last_used)
```

Parameter Description

cookie	This value is returned by dma_async_memcpy, and is used as an identifier by the DMA subsystem, so it can determine which DMA transaction the client is waiting on.
last	Returns last completed cookie, can be NULL.
used	Returns last issued cookie, can be NULL.

Return Values:

- **DMA_SUCCESS** - the DMA transaction associated with the given cookie is complete.
- **DMA_IN_PROGRESS** - the DMA transaction is not complete.
- **DMA_ERROR** - the DMA transaction failed for some reason. All in-flight transactions should be resubmitted.

dma_async_memcpy_complete

This function is used to determine if a submitted DMA request has been completed yet by the DMA engine. Typical operation would be to submit a series of memory copy operations, and then poll this function with the last cookie value returned, in order to know when it is safe to proceed.

enum dma_status_t

```
dma_memcpy_complete(struct dma_chan *chan, dma_
    cookie_t cookie, dma_cookie_t *last, dma_cookie_t *used)
```

Parameter Description

chan	An allocated dma channel pointer.
cookie	This value is returned by dma_async_memcpy, and is used as an identifier by the DMA subsystem, so it can determine which DMA transaction the client is waiting on.
last	Returns last completed cookie, can be NULL.
used	Returns last issued cookie, can be NULL.

Return Values:

- **DMA_SUCCESS** - the DMA transaction associated with the given cookie is complete.
- **DMA_IN_PROGRESS** - the DMA transaction is not complete.
- **DMA_ERROR** - the DMA transaction failed for some reason. All in-flight transactions should be resubmitted.

dma_sync_wait

This function waits until a specified DMA memory copy completes.

enum dma_status

dma_sync_wait(struct dma_chan *chan, dma_cookie_t cookie)

Parameter Description	
chan	An allocated dma channel pointer.
cookie	This value is returned by dma_async_memcpy, and is used as an identifier by the DMA subsystem, so it can determine which DMA transaction the client is waiting on.

Return Values:

- **DMA_SUCCESS** – the DMA transaction associated with the given cookie is complete.
- **DMA_IN_PROGRESS** – the DMA transaction is not complete.
- **DMA_ERROR** – the DMA transaction failed for some reason. All in-flight transactions should be resubmitted.

6.3 Client Data Structures

Client code receives pointers to dma_client and dma_chan structures. However, client code should never need to examine or change the members of these structures. The information described below is only used by IOATDMA module.

struct dma_client

This structure is used by the DMA subsystem to keep track of DMA resources allocated to the client.

```
struct dma_client {
    dma_event_callback    event_callback;
    dma_cap_mask_t        cap_mask;
    struct list_head      global_node;
};
```

Field	Description
event_callback	This is a pointer to the client's event callback, which the client passed to the DMA subsystem in dma_async_client_register.
cap_mask	This is DMA client capabilities mask.
global_node	This is used by the DMA subsystem to keep a list of all registered clients.

struct dma_chan

This structure represents a DMA channel. It is used by the DMA subsystem to keep track of which device supplied the channel, as well as what client it is allocated to, if any. It also keeps track of some per-channel statistics, and data members needed for the sysfs interface.

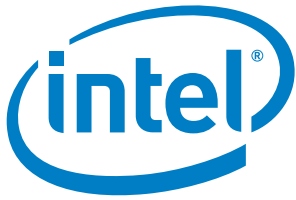
```
struct dma_chan
{
    struct dma_device *device;
    dma_cookie_t cookie;

    /* sysfs */
    int chan_id;
    struct class_device class_dev;

    struct kref refcount;
    int slow_ref;
    struct rcu_head rcu;

    struct list_head device_node;
    struct dma_chan_percpu *local;
};
```

Field	Description
device	Initialized by the driver, this is a pointer to the device this channel resides on.
cookie	Last cookie value returned to client.
chan_id	Each channel is given a number by the DMA subsystem, so they can be differentiated in sysfs.
class_dev	Required for sysfs usage.
refcount	kref, used in "bigref" slow-mode.
slow_ref	Indicates that the DMA channel is free.
rcu	The DMA channel's RCU head.
device_node	Initialized by the driver, this is used to link all channels a device has together. See struct dma_device.channels.
local	per-cpu pointer to a struct dma_chan_percpu.



www.intel.com

7. Intel QuickData Technology sysfs Interface

When the Intel QuickData Technology driver is properly loaded, there will be directories created in sysfs, under /sys/class/dma, named dma0chanX, where X is 0-3.

Channel Entries:

- in_use

1 If the DMA channel is allocated to a client, such as the network stack:

- bytes_transferred

The total number of bytes transferred by the DMA engine:

- memcpy_count

The total number of copy operations initiated.

8. Support

For general information, go to the Intel support Web site at:

<http://support.intel.com>

If an issue is identified with the released source code on the supported kernel with a supported adapter, e-mail the specific information related to the issue to linux.nics@intel.com.

9. References

www.intel.com/platforms/technologies/quickdata.htm

www.intel.com/technology/ioacceleration/index.htm

INFORMATION IN THIS DOCUMENT IS PROVIDED IN CONNECTION WITH INTEL® PRODUCTS. NO LICENSE, EXPRESS OR IMPLIED, BY ESTOPPEL OR OTHERWISE, TO ANY INTELLECTUAL PROPERTY RIGHTS IS GRANTED BY THIS DOCUMENT. EXCEPT AS PROVIDED IN INTEL'S TERMS AND CONDITIONS OF SALE FOR SUCH PRODUCTS, INTEL ASSUMES NO LIABILITY WHATSOEVER, AND INTEL DISCLAIMS ANY EXPRESS OR IMPLIED WARRANTY, RELATING TO SALE AND/OR USE OF INTEL PRODUCTS INCLUDING LIABILITY OR WARRANTIES RELATING TO FITNESS FOR A PARTICULAR PURPOSE, MERCHANTABILITY, OR INFRINGEMENT OF ANY PATENT, COPYRIGHT OR OTHER INTELLECTUAL PROPERTY RIGHT. Intel products are not intended for use in medical, life saving, life sustaining, critical control or safety systems, or in nuclear facility applications.

Intel may make changes to specifications and product descriptions at any time, without notice. Designers must not rely on the absence or characteristics of any features or instructions marked "reserved" or "undefined." Intel reserves these for future definition and shall have no responsibility whatsoever for conflicts or incompatibilities arising from future changes to them. The information here is subject to change without notice. Do not finalize a design with this information.

This specification, as well as the software described in it, is furnished under license and may only be used or copied in accordance with the terms of the license. The information in this document is furnished for informational use only, is subject to change without notice, and should not be construed as a commitment by Intel Corporation. Intel Corporation assumes no responsibility or liability for any errors or inaccuracies that may appear in this document or any software that may be provided in association with this document.

Intel processor numbers are not a measure of performance. Processor numbers differentiate features within each processor family, not across different processor families. See www.intel.com/products/processor_number for details.

The Intel processor/chipset families may contain design defects or errors known as errata, which may cause the product to deviate from published specifications. Current characterized errata are available on request.

Copies of documents, which have an order number and are referenced in this document, or other Intel literature, may be obtained by calling 1-800-548-4725, or by visiting Intel's Web Site.

Copyright © 2008 Intel Corporation. All rights reserved. Intel, the Intel logo, the Intel Leap ahead logo, and Intel Xeon are trademarks of Intel Corporation in the U.S. and other countries.

*Other names and brands may be claimed as the property of others. Printed in the United States. 0608/DU/HBD/PDF 320045-001US